

Problem 1

Problem Description

In this problem you will use PCA and TSNE to apply dimensionality reduction to 64x64 images of signed distance fields (SDFs) on parts belonging to 8 different classes. Each class is topologically similar, with some variation in void size and shape. These signed distance fields are helpful in the prediction of internal stress fields in the parts. You will also apply KNN to predict the class of the part with the reduced space.

Fill out the notebook as instructed, making the requested plots and printing necessary values.

You are welcome to use any of the code provided in the lecture activities.

Summary of deliverables:

- 3x8 subplot visualization of the first 3 samples from each of the 8 classes
- Bar plot of the variance explained for the first 25 PCs and the number of PCs required to explain > 90% of the variance in the training data
- 4x8 subplot visualization of reconstructed samples using 3, 10, 50 and all PCs on the first sample from each of the 8 classes in the test set
- Test accuracy of KNN classifier trained on the 3D, 10D, and 50D PCA reduced feature spaces
- Plot of the 2D TSNE reduced feature space
- Test accuracy of the KNN classifier trained on the 2D TSNE reduced feature space
- Discussion questions 1 and 2

Imports and Utility Functions:

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from scipy import io

from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
from sklearn.neighbors import KNeighborsClassifier
```

```

from sklearn.model_selection import train_test_split

def dataLoader(filepath):
    # Load and flatten the SDF dataset
    mat = io.loadmat(filepath)
    data = []
    for i in range(800):
        sdf = mat["sdf"][i][0].T
        data.append(sdf.flatten())
    data = np.vstack(data)
    # Assign Labels
    labels = np.repeat(np.arange(8), 100)
    return data, labels

def plot_sdf(data, ax = None, title = None):
    # If no axes, make them
    if ax is None:
        ax = plt.gca()
    # Reshape image data into square
    sdf = data.reshape(64,64)
    # Plot image, with bounds of the SDF values for the entire dataset
    ax.imshow(sdf, vmin=-0.31857, vmax=0.206349, cmap="jet")
    ax.axis('off')
    # If there is a title, add it
    if title:
        ax.set_title(title)

```

Visualization

Using the provided `dataLoader()` function, load the data and labels from `sdf_images.mat`. The returned data will contain 800 samples, with 4096 features. Then, using the provided `plot_sdf()` function, generate a 3x8 subplot figure containing visualizations of the first 3 SDFs in each class.

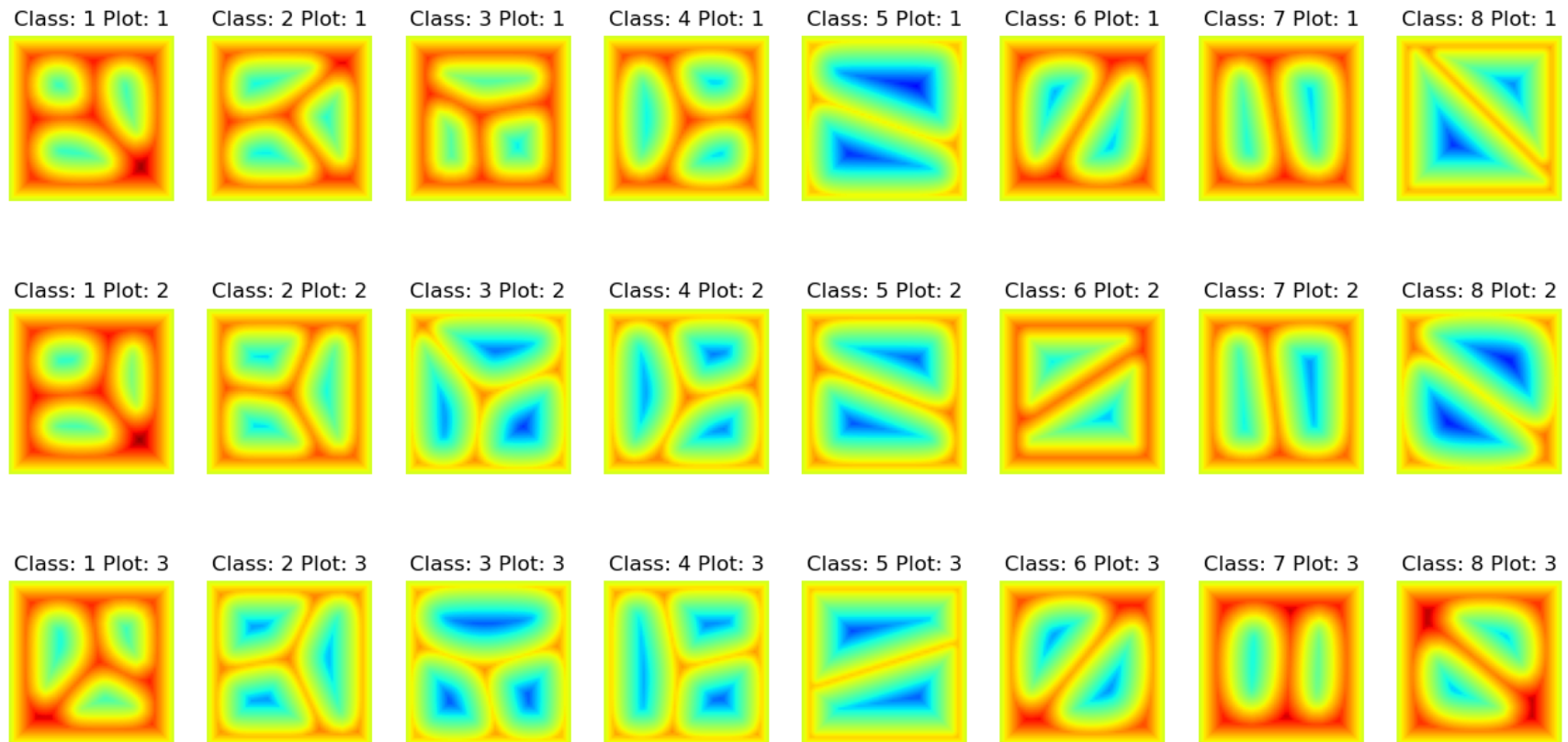
```

In [2]: # YOUR CODE GOES HERE
data, labels = dataLoader("data/sdf_images.mat")

fig, ax = plt.subplots(3, 8, figsize=(16, 8))

for i in range(8):
    for j in range(3):
        index = 100*i + j
        plot_sdf(data[index], ax[j,i], title = f'Class: {i+1} Plot: {j+1}')

```



Explained Variance

Use `train_test_split()` to partition the data and labels into a training and test set with `test_size = 0.2` and `random_state = 0`. Then train a PCA model on the training data and generate a bar plot of the variance explained for the first 25 principal components. Determine the number of principal components required to explain > 90% of the variance in the training data.

```
In [3]: # YOUR CODE GOES HERE
X_train,X_test,y_train,y_test = train_test_split(data,labels,test_size = 0.2,random_state = 0)

pca = PCA()
pca.fit(X_train)
X_train_pca = pca.transform(X_train)
```

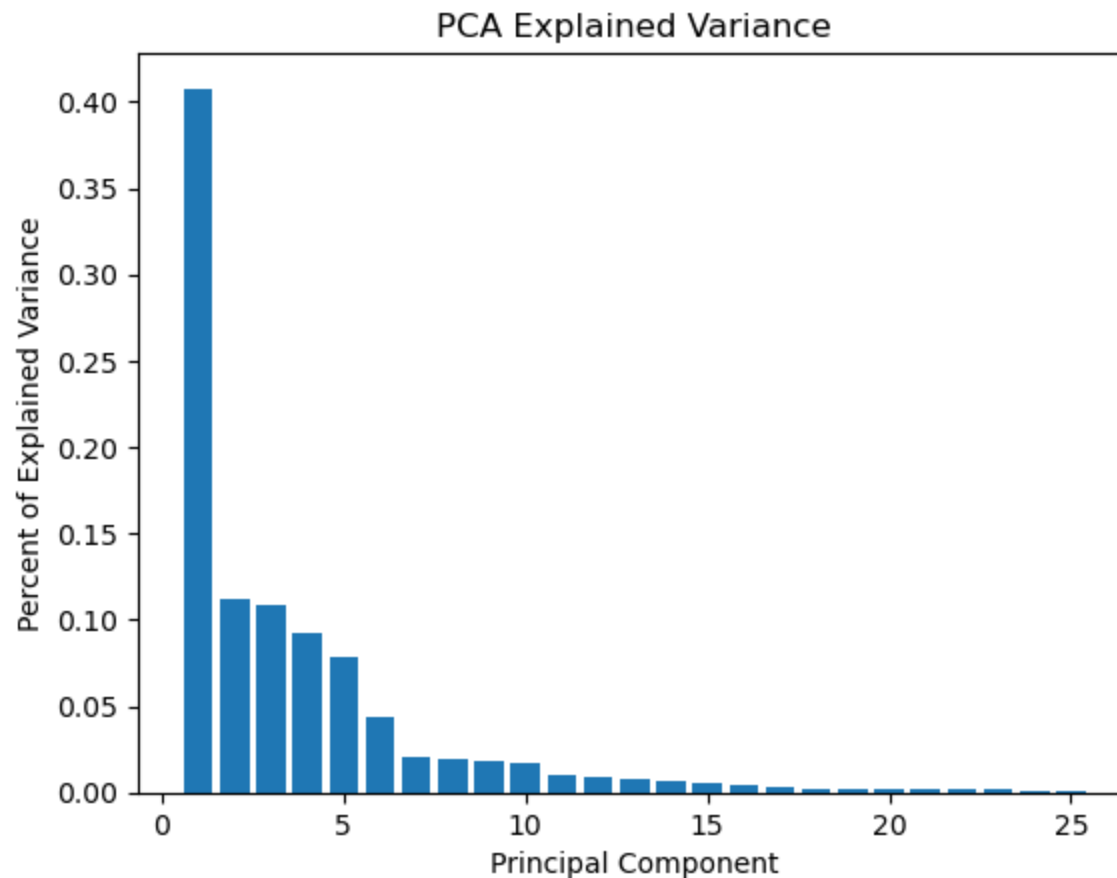
```
plt.figure()
plt.bar(range(1,26),pca.explained_variance_ratio_[:25])
plt.xlabel("Principal Component")
plt.ylabel("Percent of Explained Variance")
plt.title("PCA Explained Variance")
plt.show()

cumulative_variance_ratio = np.cumsum(pca.explained_variance_ratio_)

num_components = 0

for i, ratio in enumerate(cumulative_variance_ratio):
    if ratio > 0.9:
        num_components = i + 1
        break

print(f"Number of principal components required to explain > 90% of the variance: {num_components}")
```



Number of principal components required to explain > 90% of the variance: 9

PCA Reconstruction

Using the training data, generate 4 PCA models using 3, 10, 50, and all of the principal components. Use these models to transform the test data into the reduced space, and then reconstruct the data from the reduced space. Plot the reconstruction for each model, on the first occurrence of each class in the test set. Your generated plot should be a 4x8 subplot figure, with each subplot title containing the class and the number of PCs used.

```
In [4]: # YOUR CODE GOES HERE

n = min(X_train.shape[0],X_train.shape[1])
num_pcs = [3,10,50,n]
```

```
fig, ax = plt.subplots(4, 8, figsize = (16, 8))

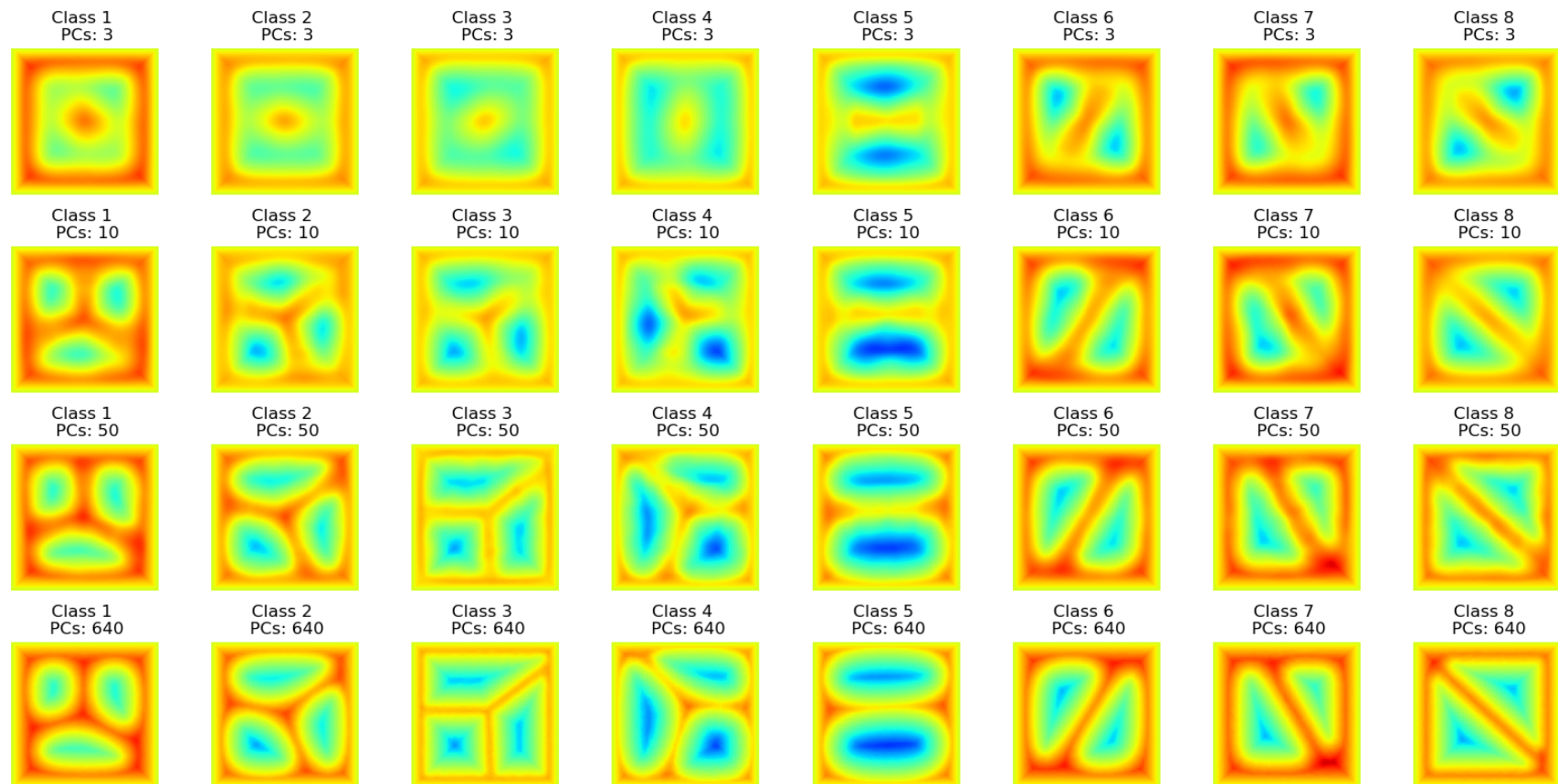
for i, num in enumerate(num_pcs):
    pca_model = PCA(n_components = num)
    pca_model.fit(X_train)

    X_train_pca = pca_model.transform(X_train)
    X_test_pca = pca_model.transform(X_test)

    X_train_recon = pca_model.inverse_transform(X_train_pca)
    X_test_recon = pca_model.inverse_transform(X_test_pca)

    for j in range(8):
        index = np.argmax(y_test == j)
        plot_sdf(X_test_recon[index], ax=ax[i, j], title = f'Class {j+1} \n PCs: {num}')

plt.tight_layout()
plt.show()
```



KNN on PCA Reduced Data

Now train a KNN classifier to predict the class of the 3D, 10D, and 50D PCA reduced data. You should train the KNN on the reduced training data, and report the prediction accuracy on the test set. You will also need to determine the `n_neighbors` parameter for your KNN classifier that gives good results.

```
In [5]: # YOUR CODE GOES HERE
from sklearn.metrics import accuracy_score

for num in [3, 10, 50]:

    pca = PCA(n_components=num)
    X_train_pca = pca.fit_transform(X_train)
    X_test_pca = pca.transform(X_test)
```

```

final = 0
n_neighbor = 0
for k in range(2,11):

    knn = KNeighborsClassifier(n_neighbors = k)
    knn.fit(X_train_pca,y_train)
    y_pred = knn.predict(X_test_pca)
    accuracy = accuracy_score(y_test,y_pred)*100
    if accuracy > final:
        final = accuracy
        n_neighbor = k

print(f'Accuracy for {num}D for test_data with n_neighbors={n_neighbor} is {final}% ')

```

Accuracy for 3D for test_data with n_neighbors=7 is 71.25%
 Accuracy for 10D for test_data with n_neighbors=8 is 91.25%
 Accuracy for 50D for test_data with n_neighbors=8 is 91.25%

TSNE Visualization

First reduced the full dataset to 50D using PCA, and then further reduced the data to 2D using TSNE. Plot the 2D reduced feature space with a scatter plot, coloring each point according to its class.

```

In [6]: # YOUR CODE GOES HERE

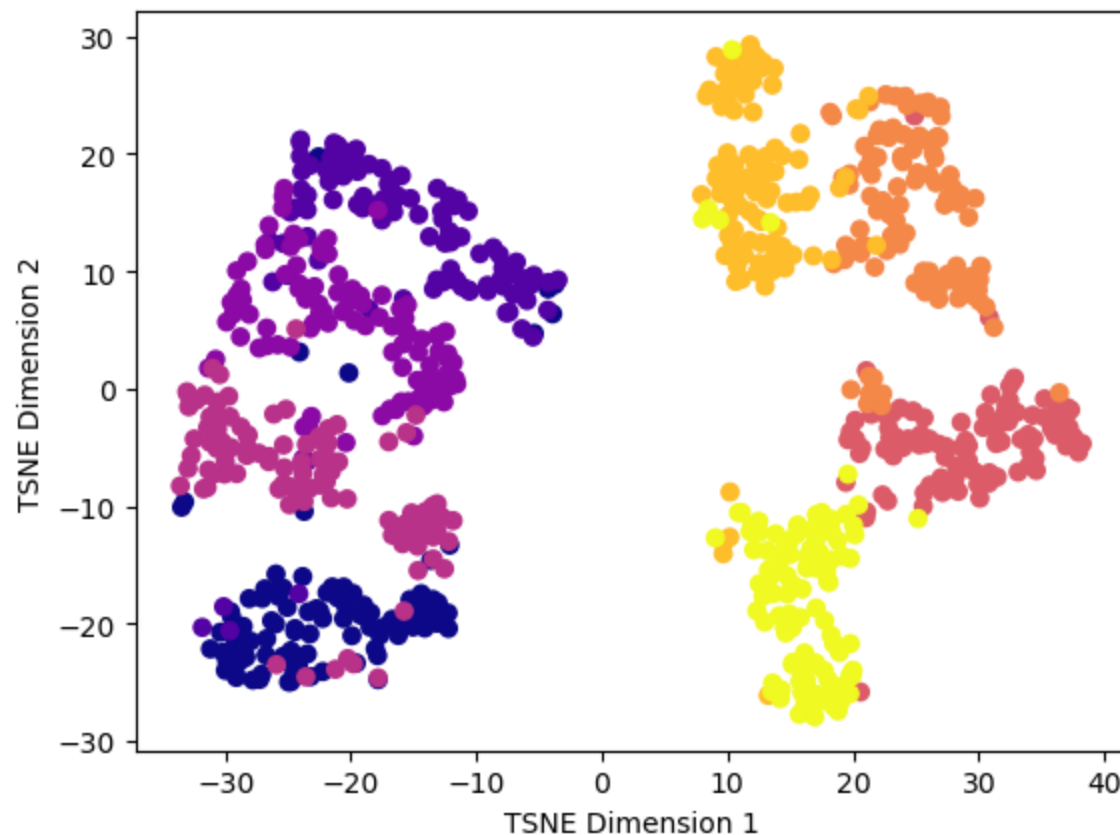
pca = PCA(n_components = 50)
pca.fit(data,labels)
data_50 = pca.transform(data)

tsne = TSNE(n_components = 2)
data_tsne = tsne.fit_transform(data_50)

plt.figure()
plt.scatter(data_tsne[:,0],data_tsne[:,1],c=labels,cmap = 'plasma')
plt.xlabel("TSNE Dimension 1")
plt.ylabel("TSNE Dimension 2")

plt.show()

```

KNN on PCA/TSNE Reduced Data

Using the same 2D PCA/TSNE data, split the data into train and test data and labels using `train_test_split` with a `random_state = 0` parameter so you have the same train/test partition as before. Then, train a KNN on this 2D feature space with the training set, and report the KNN classifier accuracy on the test set. Again, you will need to determine the `n_neighbors` parameter in the KNN classifier that gives good results.

```
In [7]: # YOUR CODE GOES HERE
X_train,X_test,y_train,y_test = train_test_split(data_tsne,labels,test_size = 0.2,random_state = 0)

final = 0
n_neighbor = 0
for k in range(2,11):
```

```
knn = KNeighborsClassifier(n_neighbors = k)
knn.fit(X_train,y_train)
y_pred = knn.predict(X_test)
accuracy = accuracy_score(y_test,y_pred)*100
if accuracy > final:
    final = accuracy
    n_neighbor = k

print(f'Accuracy on the test set for n_neighbors = {n_neighbor} is {final}%')
```

Accuracy on the test set for n_neighbors = 6 is 90.0%

Discussion

1. Discuss how the number of principal components relates to the quality of reconstruction of the data. Using all of the principal components, should there be any error in the reconstruction of a sample from the training data? What about in the reconstruction of an unseen sample from the testing data?
2. Discuss how you determined `k`, the number of neighbors in your KNN models. Why do we perform dimensionality reduction to our data before feeding it to our KNN classifier?

Your response goes here

1. As the number of principal components increase, the quality of reconstruction data improves. If all the principal components are used then there will be no error in the reconstruction of a sample from training data. This means that it will be lossless. This also applies to the unseen sample from the testing data, it will be lossless.
2. The `k` i.e., the number of neighbors in KNN models was selected by looping through a range of `k`-values and seeing which one gives the best performance. This helped in providing an accuracy around 90% which means that the best performance is seen around 90%. Not a higher value of `k` was taken because it can lead to underfitting of the data. Before doing the KNN classification, the dimensionality reduction is done because KNN is sensitive to dimensionality. As the number of features or dimensions increases it can downgrade the performance of the model. As the dimension increases, the distance between the data points becomes less meaningful which can lead to computational complexity and overfitting of the data. Along with that, PCA helps in capturing the most important features and therefore the generalization of data becomes better which gives us a good fit.

In []: